

# *E Pluribus Unum Ex Machina* Learning from Many Collider Events at Once

~ work with Jesse Thaler in PRD 103 (2021) 116013, 2101.07263 ~

## Benjamin Nachman

*Lawrence Berkeley National Laboratory*

[bpnachman.com](http://bpnachman.com)

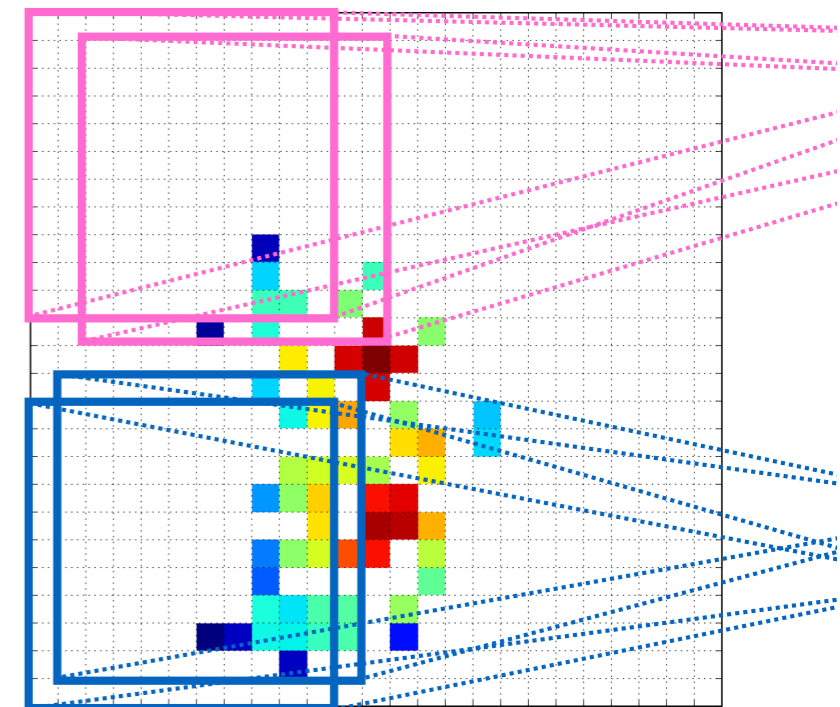
[bpnachman@lbl.gov](mailto:bpnachman@lbl.gov)



@bpnachman



bnachman



CMS Machine  
Learning Town Hall

July 27, 2021

- The core idea of this paper: *you can't (formally) gain by learning with multiple events at once,  $f = f(x_1, \dots, x_N)$ .*
  - ➔ This may be obvious to some of you and it may violate the intuition of others.
  - ➔ There may be practical situations where it is better to train a classifier/regressor to process multiple events at once, especially if training can be made simpler.
- In addition to explaining the core idea, I'll describe some other statistical facts that may be useful for classification, inference, etc.

# The statistical structure of collisions



3

To a very good approximation, collisions are independent:

$$p(\{x_1, \dots, x_N\} | \theta) = \prod_{i=1}^N p(x_i | \theta)$$

$x_i$  are features to represent events  
(could be low-level or high-level)

# The statistical structure of collisions

4

This means that the **per-ensemble likelihood ratio can be written as a product of per-event likelihood ratios.**

$$\frac{p(\{x_1, \dots, x_N\} | \theta_A)}{p(\{x_1, \dots, x_N\} | \theta_B)} = \prod_{i=1}^N \frac{p(x_i | \theta_A)}{p(x_i | \theta_B)}$$

This explains the informational equivalence of an optimal per-ensemble learning (left) and a per-event learning (right)

*Collections of  $N$  events have more info than single events, but a per-event classifier applied  $N$  times has access to the same info.*

# Brief review: per-instance classification

5

We start with a loss function(al):

$$L[f] = - \sum_{i \in \text{class } A} A(f(x_i)) - \sum_{i \in \text{class } B} B(f(x_i))$$

It is often useful to consider the continuum limit:

$$L[f] = - \int dx (p(x | \theta_A) A(f(x)) + p(x | \theta_B) B(f(x)))$$

Typically,  $A = \log$  and  $B = 1 - \log$

# Brief review: per-instance classification



Functional optimization shows that

$$\frac{B'(f(x))}{A'(f(x))} = \frac{p(x | \theta_A)}{p(x | \theta_B)} \leftarrow \text{optimal classifier}$$

This means that there are many loss functionals (defined by  $A$  and  $B$ ) which result in an optimal classifier.

in particular, as long as  $B'(f)/A'(f)$  is a monotonic rescaling of  $f$  and the overall loss is convex

# Brief review: per-instance classification



$$\frac{B'(f(x))}{A'(f(x))} = \frac{p(x|\theta_A)}{p(x|\theta_B)} \quad \leftarrow \text{optimal classifier}$$

Since it simplifies much of the notation, consider the “maximum likelihood classifier” (MLC) loss\* with  $A(f) = \log(f)$  and  $B(f) = 1-f$ . Then,

$$\operatorname{argmin}_f L_{\text{MLC}}[f] = \frac{p(x|\theta_A)}{p(x|\theta_B)}$$

*The loss value itself is the KL divergence between the conditional probabilities*

\*This was introduced in its exponential form by R. T. D’Agnolo and A. Wulzer in PRD 99 (2019) 015014, 1806.02350

# Per-ensemble classification



One could insert  $N$  events instead of one into the MLC loss:

$$L_{\text{MLC}}[f_N] = - \int d^N x \left( \overline{x} \mid \theta_A \right) \log f_N(\overline{x}) + p(\overline{x} \mid \theta_B)(1 - f_N(\overline{x}))$$

Unsurprisingly,

$$\operatorname{argmin}_{f_N} L_{\text{MLC}}[f_N] = \frac{p(\overline{x} \mid \theta_A)}{p(\overline{x} \mid \theta_B)}$$



# Per-ensemble classification



$$\operatorname{argmin}_{f_N} L_{\text{MLC}}[f_N] = \frac{p(\bar{x} | \theta_A)}{p(\bar{x} | \theta_B)}$$

However, we can also build this up from one event:

$$f_{1 \rightarrow N} \equiv \prod_{i=1}^N f_1(x_i) \rightarrow \frac{p(\bar{x} | \theta_A)}{p(\bar{x} | \theta_B)}$$

If  $x$  is  $k$ -dimensional, then  $f_1$  only requires samples of **dimension  $k$**  while  $f_N$  requires samples of **dimension  $k^N$**

# Per-ensemble classification

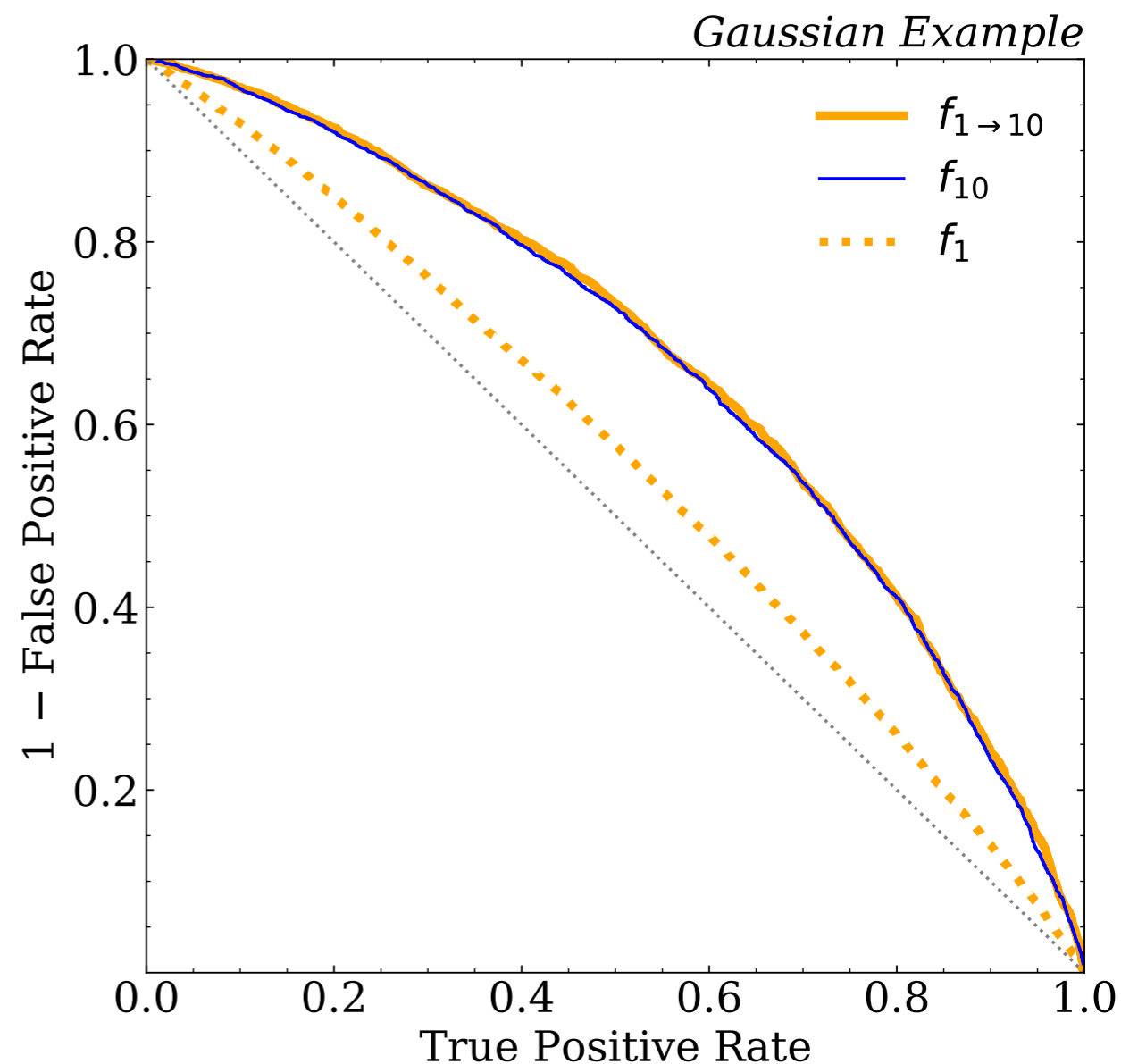
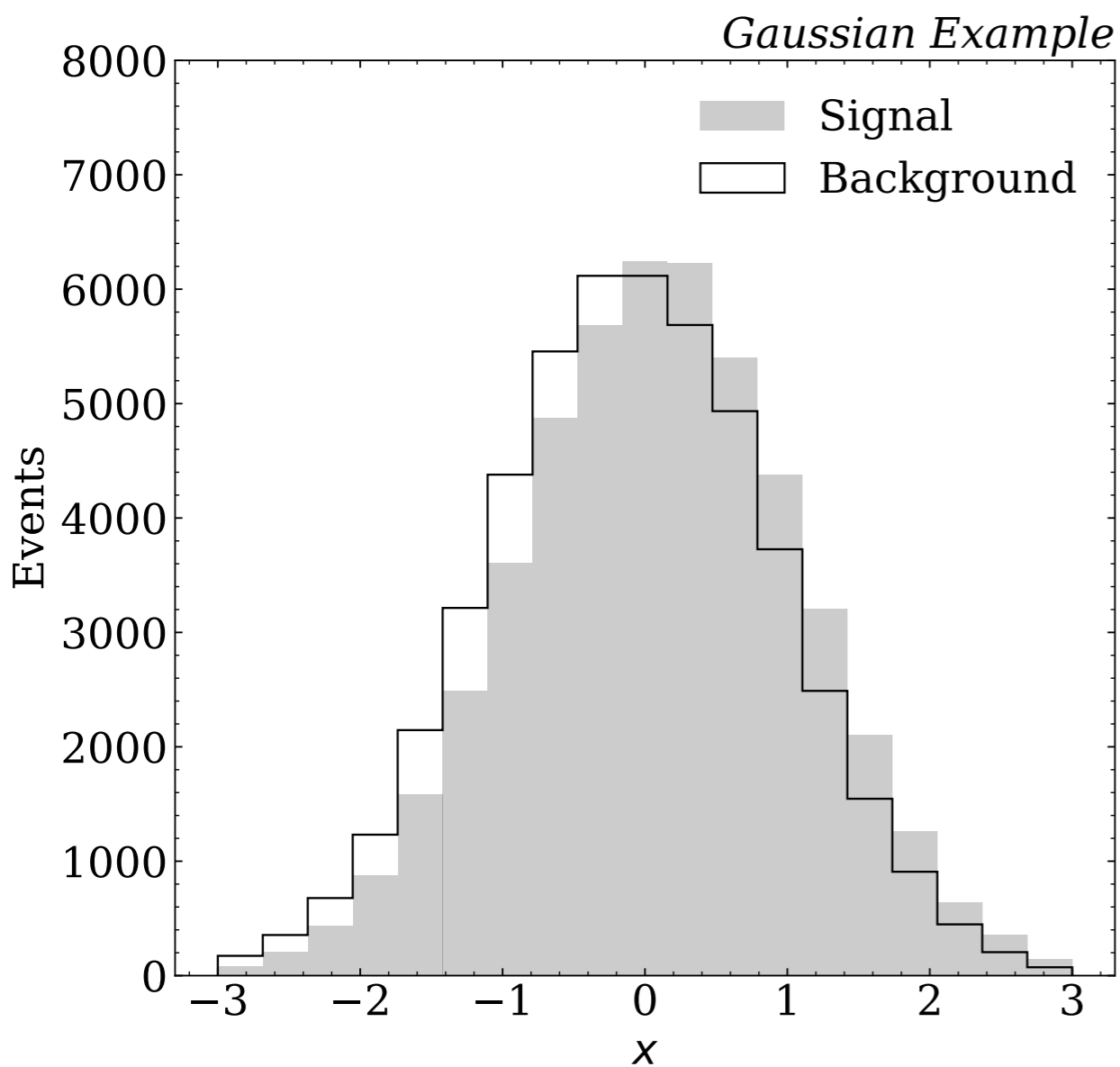
10




$$f_{1 \rightarrow N} \equiv \prod_{i=1}^N f_1(x_i) \rightarrow \frac{p(\vec{x} | \theta_A)}{p(\vec{x} | \theta_B)}$$

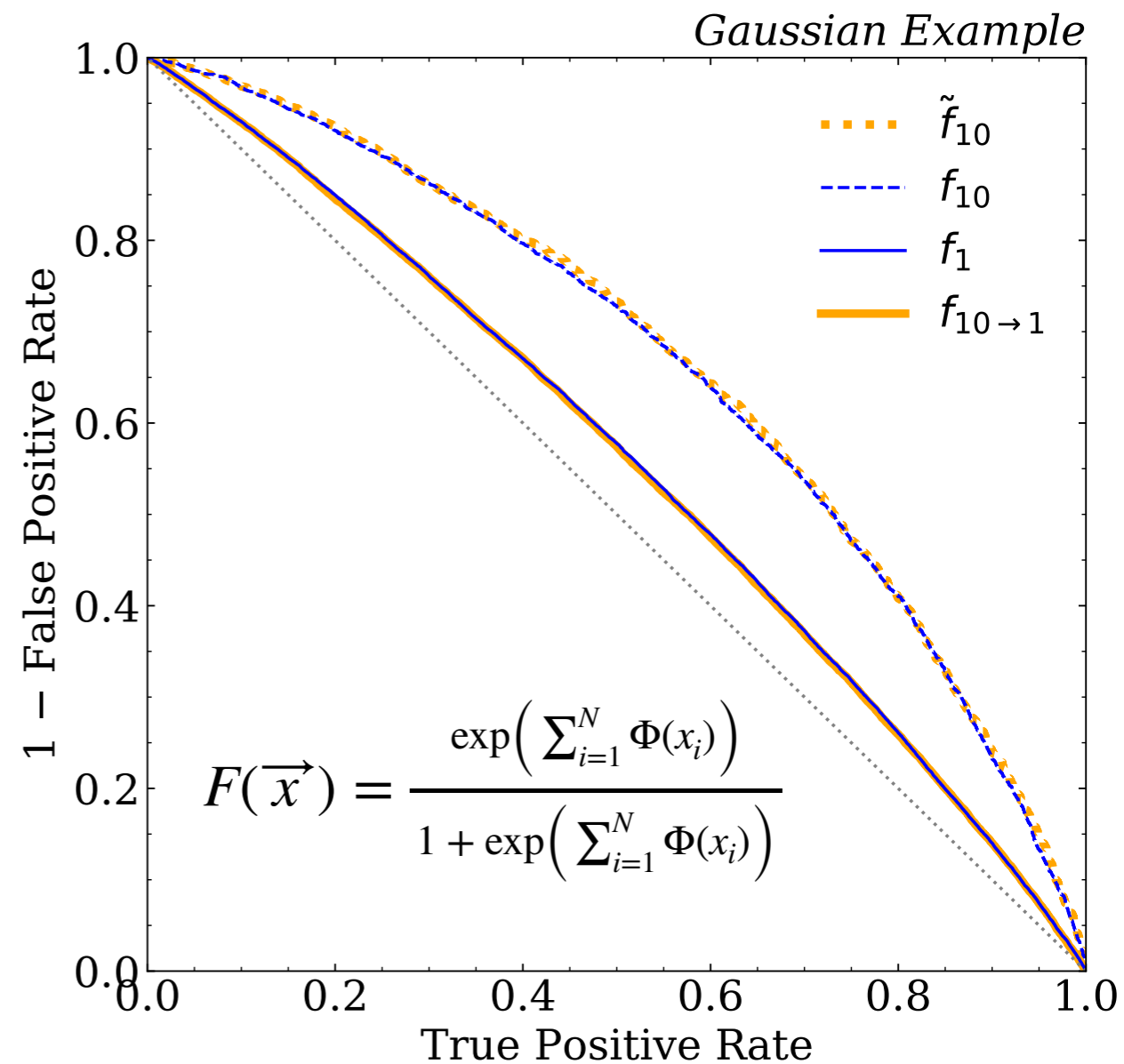
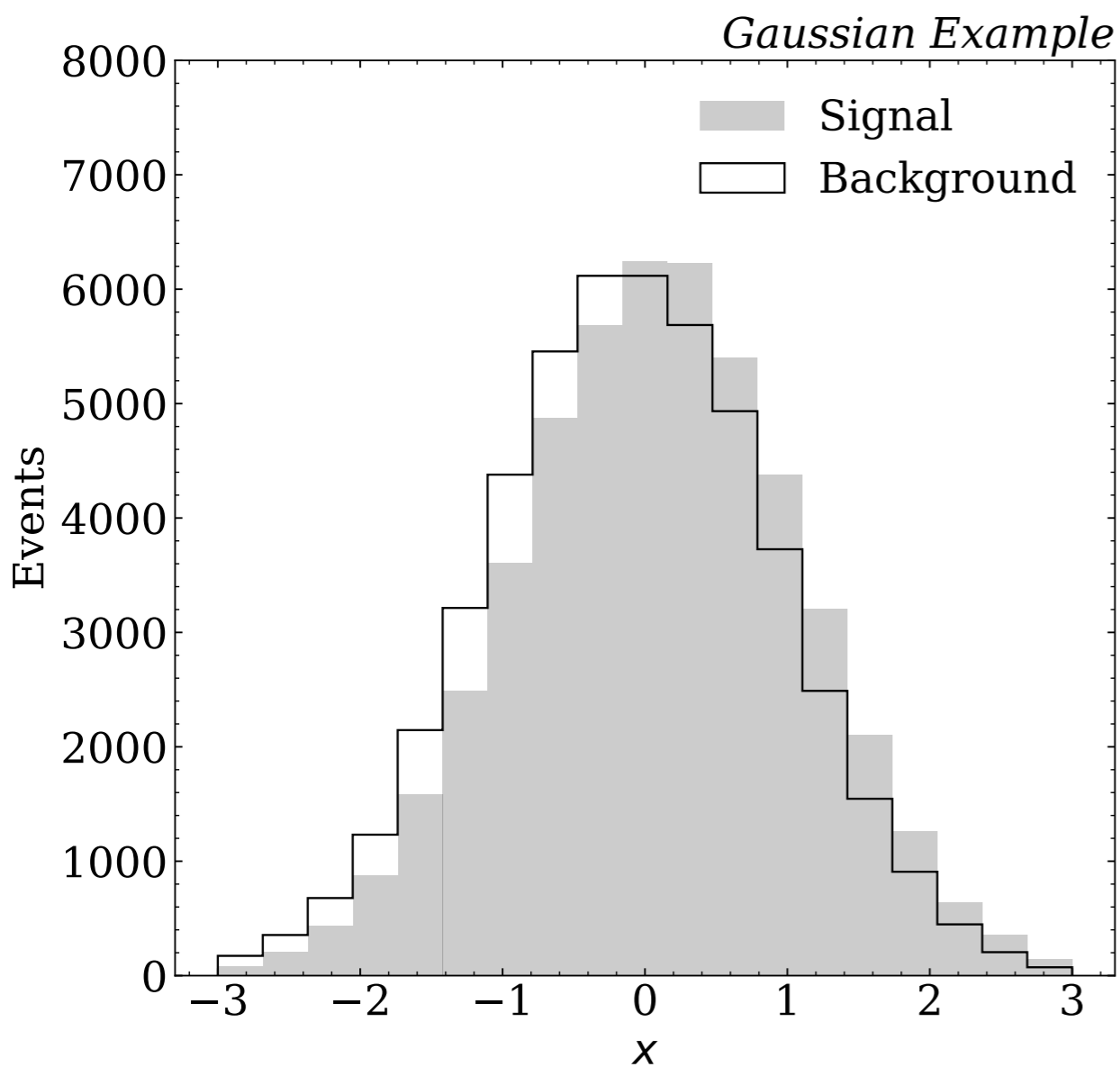
It is also possible to go in the converse direction:

$$\tilde{f}_N(\vec{x}) \equiv \prod_{i=1}^N f_{N \rightarrow 1}(x_i)$$

Where the entire righthand side is used to minimize the ensemble MLC loss. This construction is at the core of set-based neural networks like **Deep Sets**.

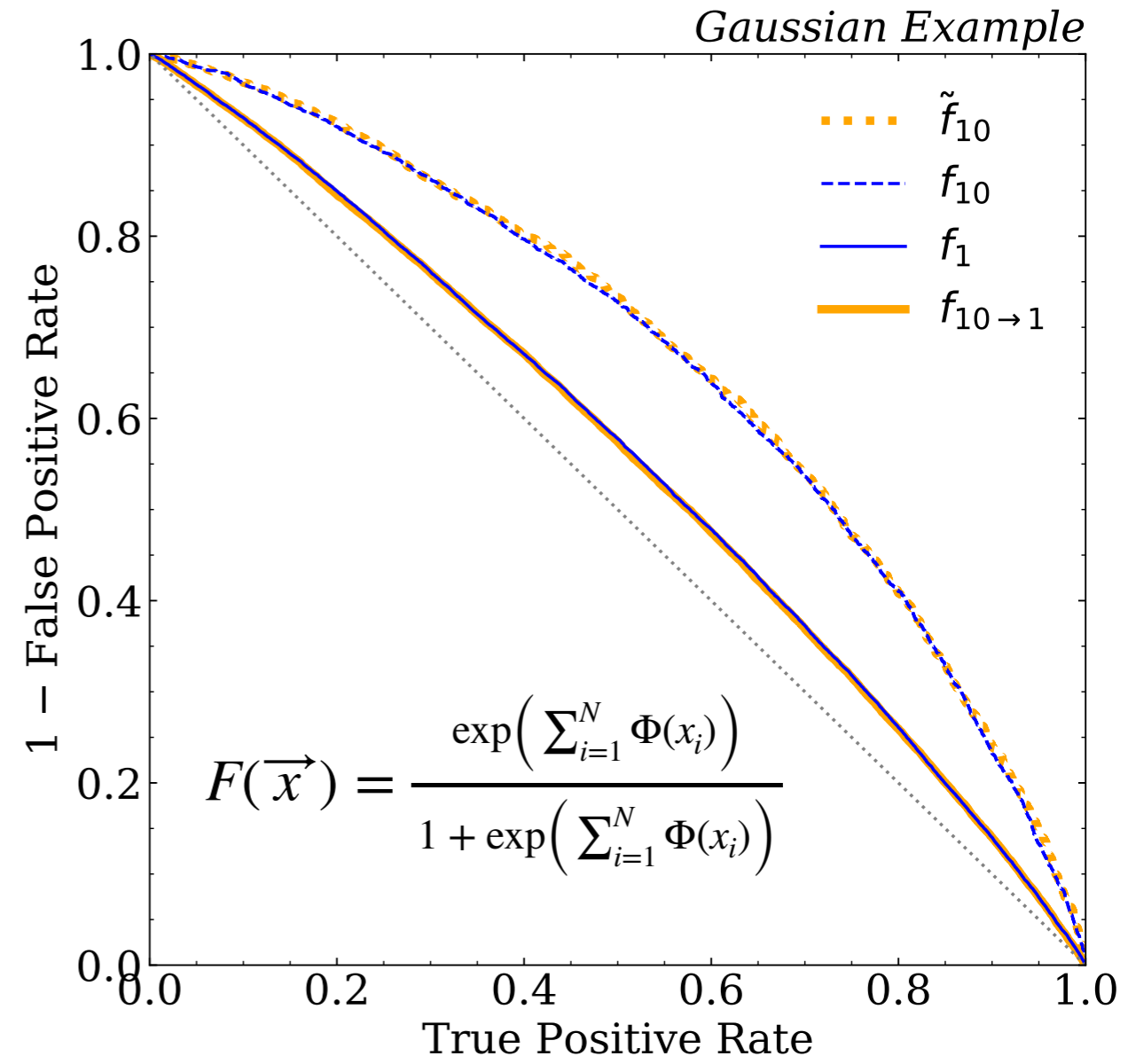
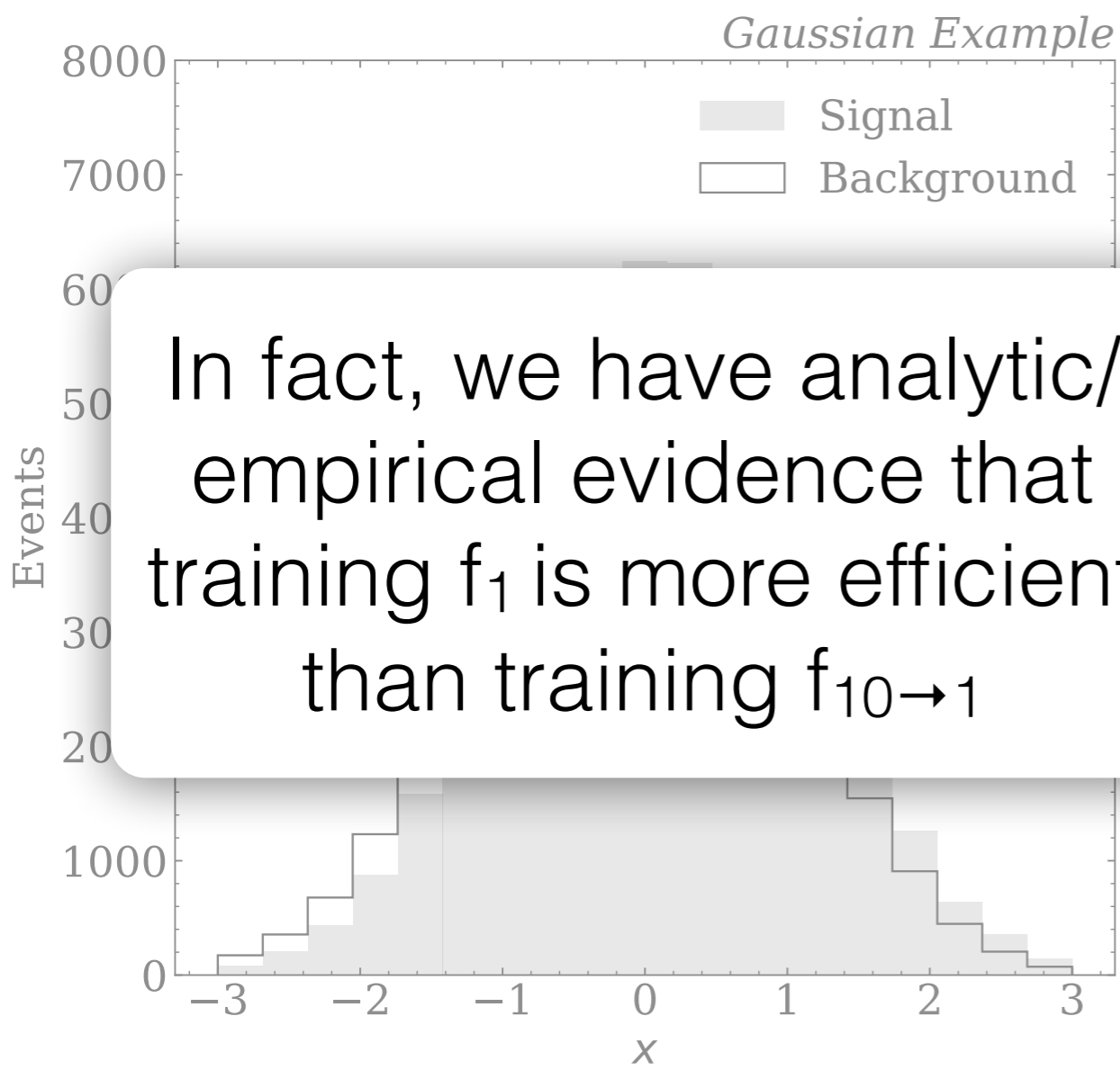


-   $f_{1 \rightarrow 10}$  per-event classifier applied 10 times
-   $f_{10}$  per-(10-event) classifier
-   $f_1$  per-event classifier



—  $f_{10 \rightarrow 1}$

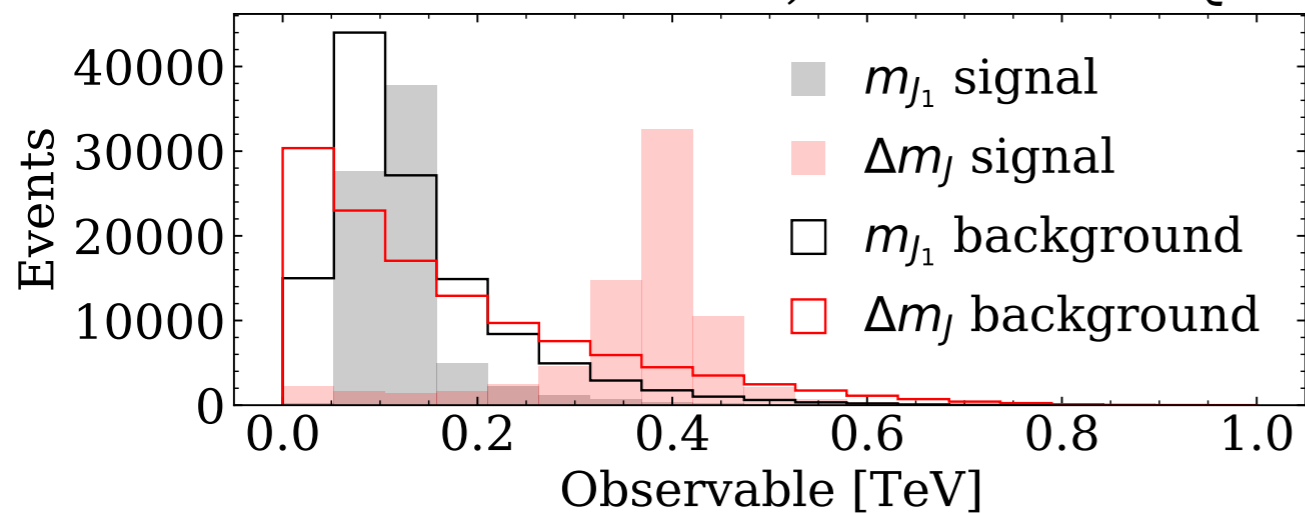
Train a **particle flow network**  $F(\Phi(x))$  with latent space dimension 1 and non-trainable  $F$



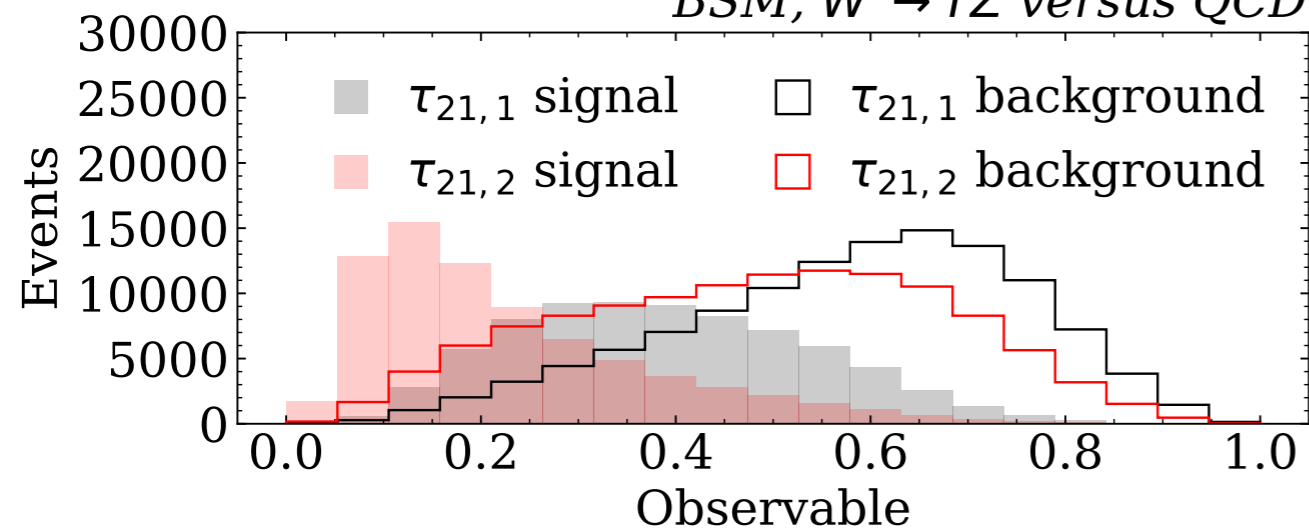
—  $f_{10 \rightarrow 1}$

Train a **particle flow network**  $F(\Phi(x))$  with latent space dimension 1 and non-trainable  $F$

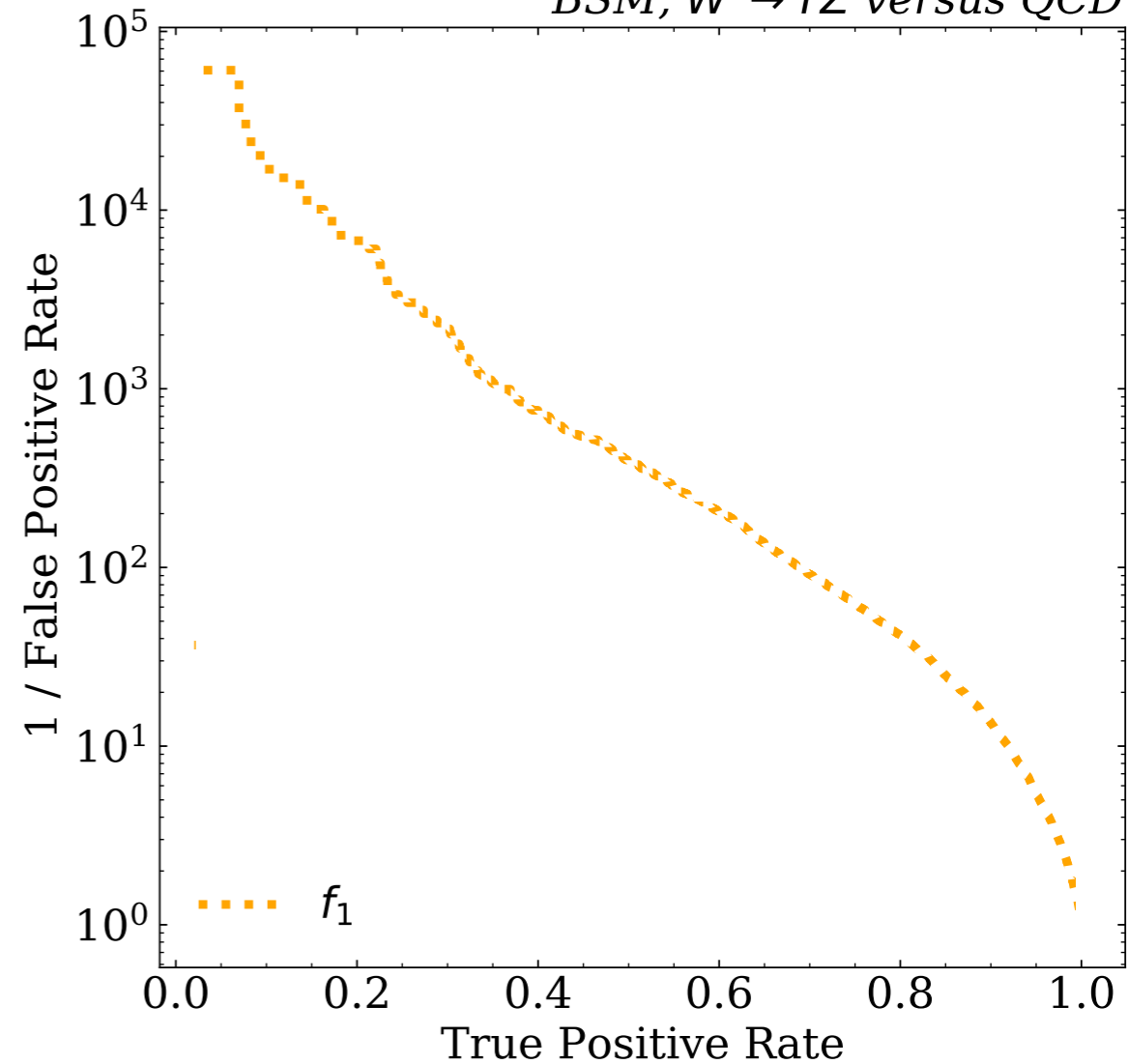
*BSM,  $W' \rightarrow YZ$  versus QCD*








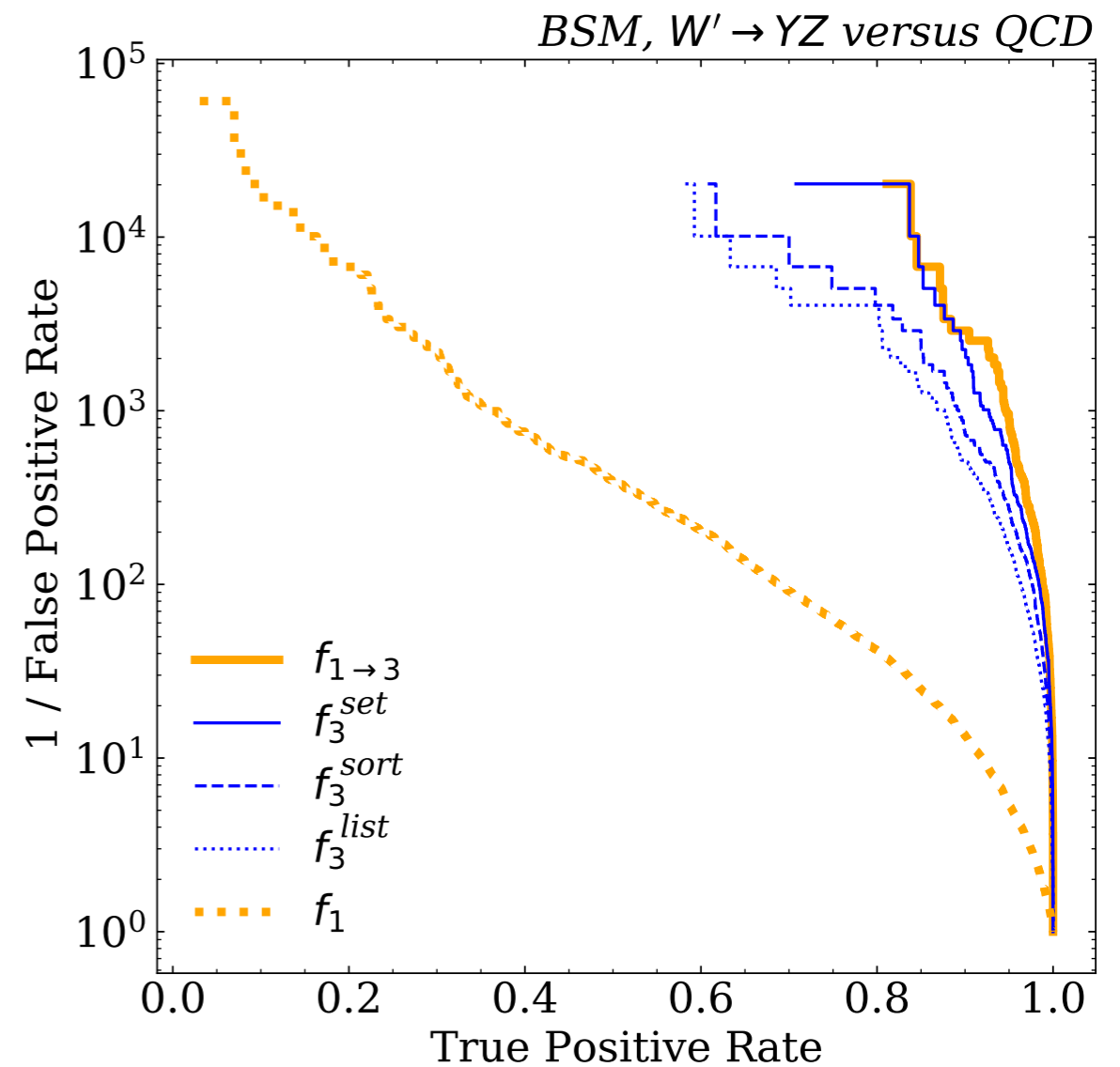
*BSM,  $W' \rightarrow YZ$  versus QCD*








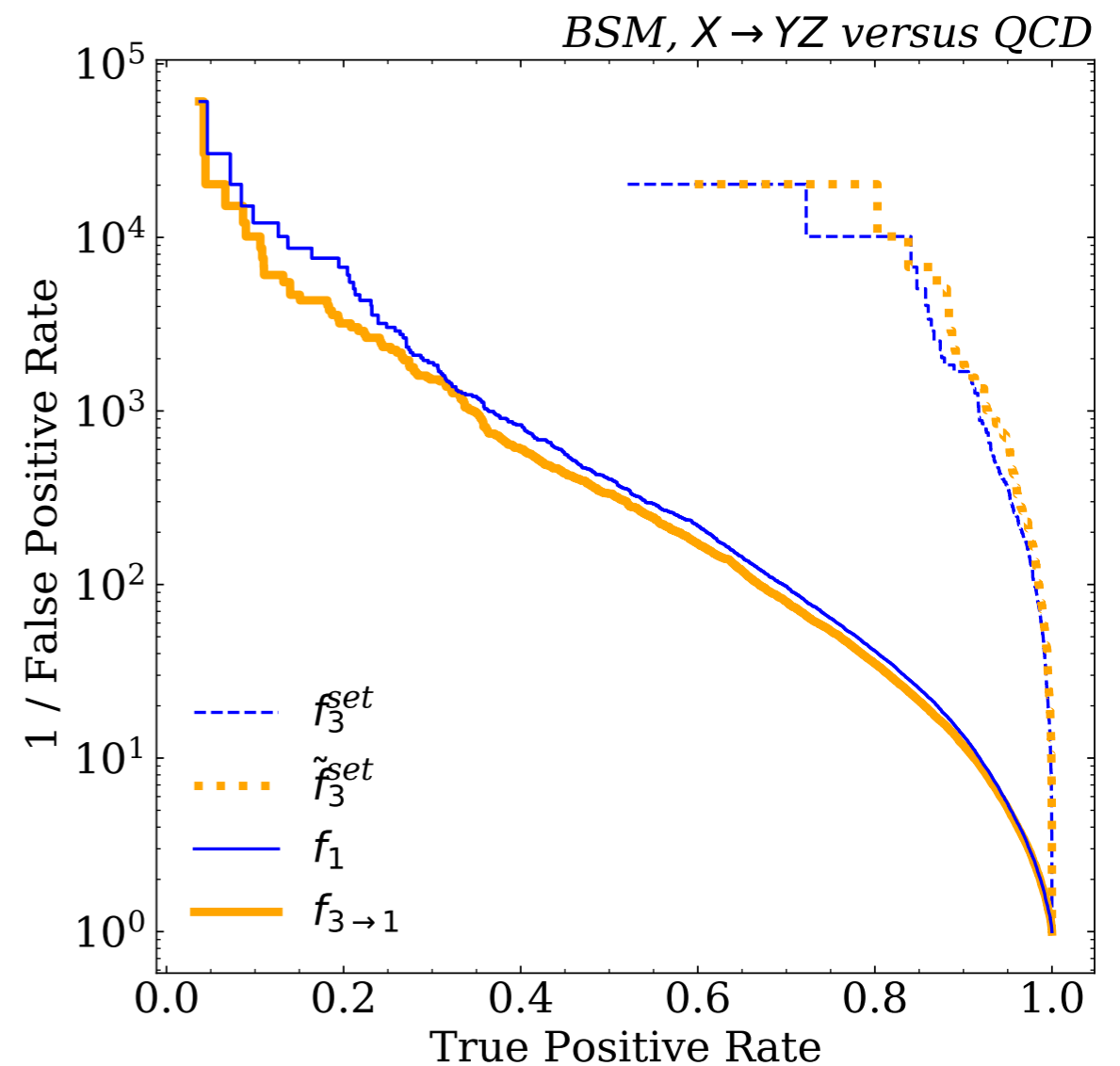
*BSM,  $W' \rightarrow YZ$  versus QCD*



-   $f_{1 \rightarrow 3}$  per-event classifier applied 3 times
-   $f_3^{set}$  per-3-event classifier, using PFN across events
-   $f_3^{sort}$  per-3-event classifier, events sorted by  $m_{J_1}$
-   $f_3^{list}$  per-3-event classifier
-   $f_1$  per-event classifier



	$f_{1 \rightarrow 3}$	per-event classifier applied 3 times
	$f_3^{set}$	per-3-event classifier, using PFN across events
	$f_3^{sort}$	per-3-event classifier, events sorted by $m_{J_1}$
	$f_3^{list}$	per-3-event classifier
	$f_1$	per-event classifier



$\tilde{f}_3^{set}$  1D latent space,  
Non-learnable F



# Per-ensemble regression I

17

Most common strategy is maximum likelihood:

$$\theta_{\text{ML}} = \operatorname{argmax}_{\theta} p(\vec{x} \mid \theta)$$

(see paper for other examples beyond regression,  
e.g. estimating the mutual information)

Most common strategy is maximum likelihood:

$$\theta_{\text{ML}} = \operatorname{argmax}_{\theta} p(\vec{x} | \theta)$$

One way of doing this would be to train a **parameterized classifier** using the MLC loss

$$f(x, \theta) = \frac{p(x | \theta)}{p(x | \theta_0)} \rightarrow \theta_{\text{ML}} = \operatorname{argmin}_{\theta} \left\{ - \sum_{i=1}^N \log f(x_i, \theta) \right\}$$

fixed reference

# Per-ensemble regression II

19

Step 1

$$f(x, \theta) = \frac{p(x | \theta)}{p(x | \theta_0)}$$

Step 2

$$\theta_{\text{ML}} = \operatorname{argmin}_{\theta} \left\{ - \sum_{i=1}^N \log f(x_i, \theta) \right\}$$

This requires two steps (amortized). Can we do it in one step?

$$\operatorname{argmax}_{\theta} \left\{ \min_f L_{\text{MLC}}[f] \right\} = \theta_{\text{ML}}$$

(see paper for derivation)

This does it all in one step, but does require a minimax optimization

You may also ask, why not directly regress  $\theta$  from  $N$  events?

$$L_{\text{MSE}}[g_N] = - \int d^N x p(\vec{x}, \theta) (g_N(\vec{x}) - \theta)^2$$

MSE = mean squared error

This is prior-dependent, but it is well-known that formally:

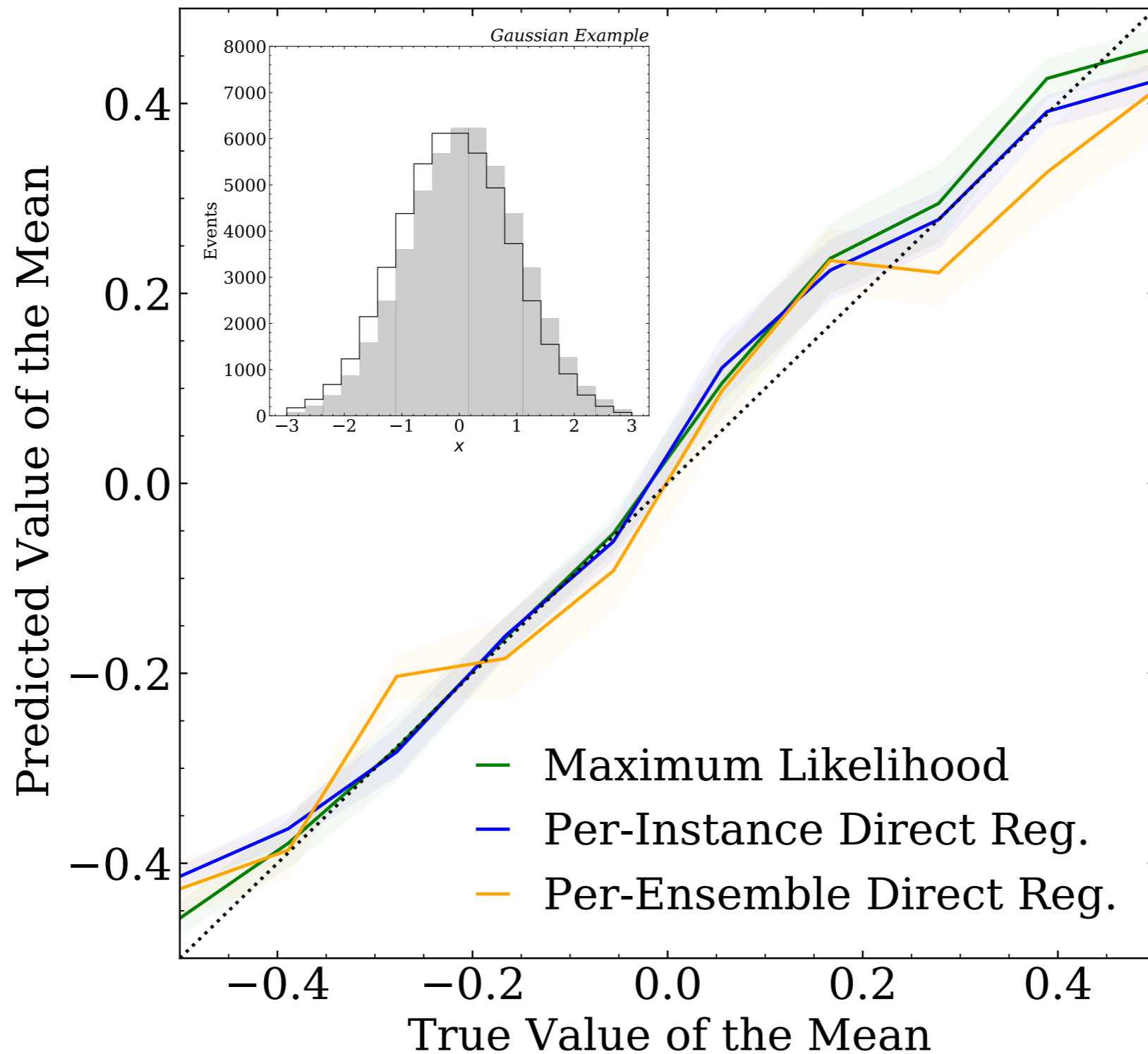
$$g_N(\vec{x}) = \langle \theta | \vec{x} \rangle$$

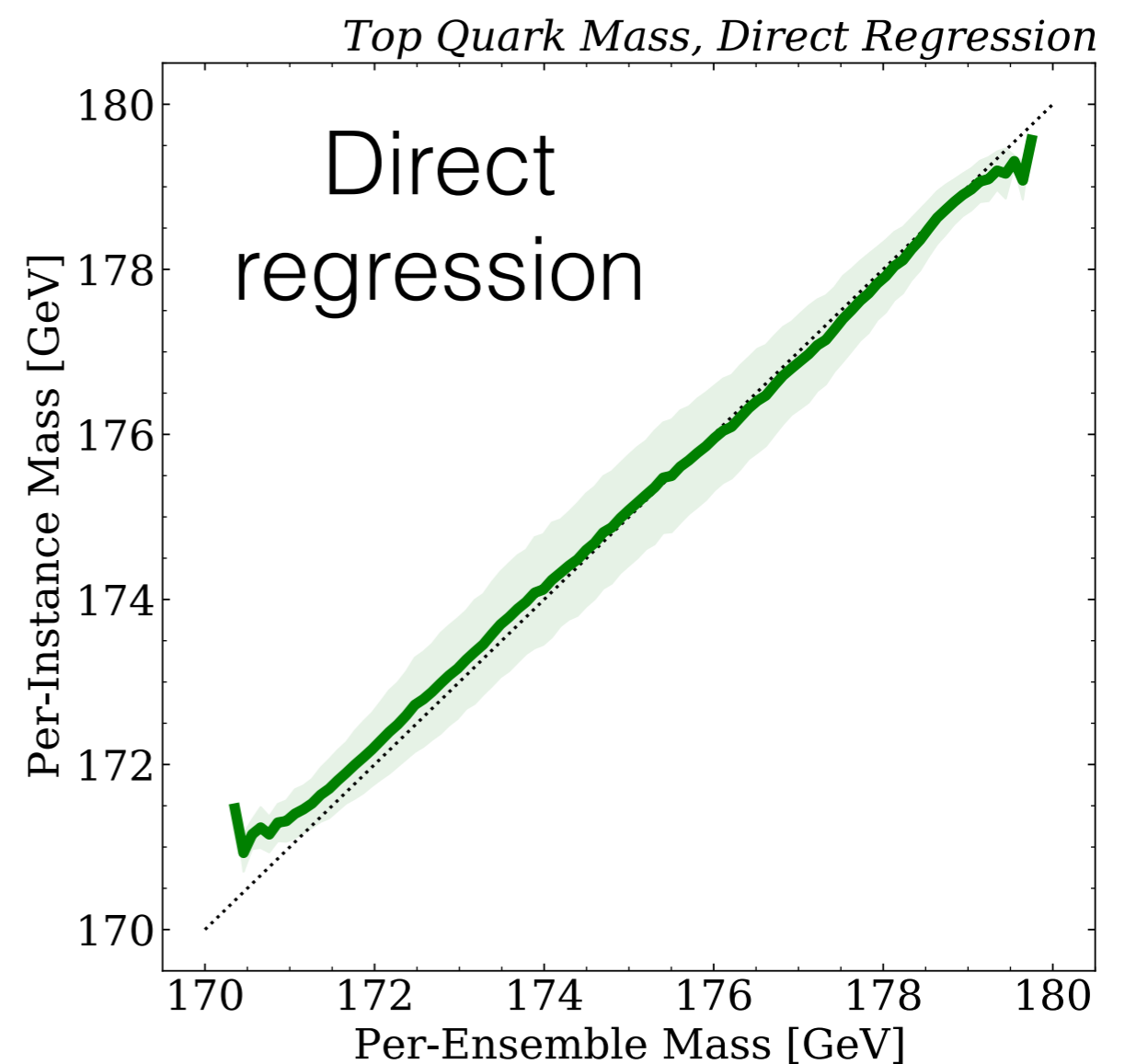
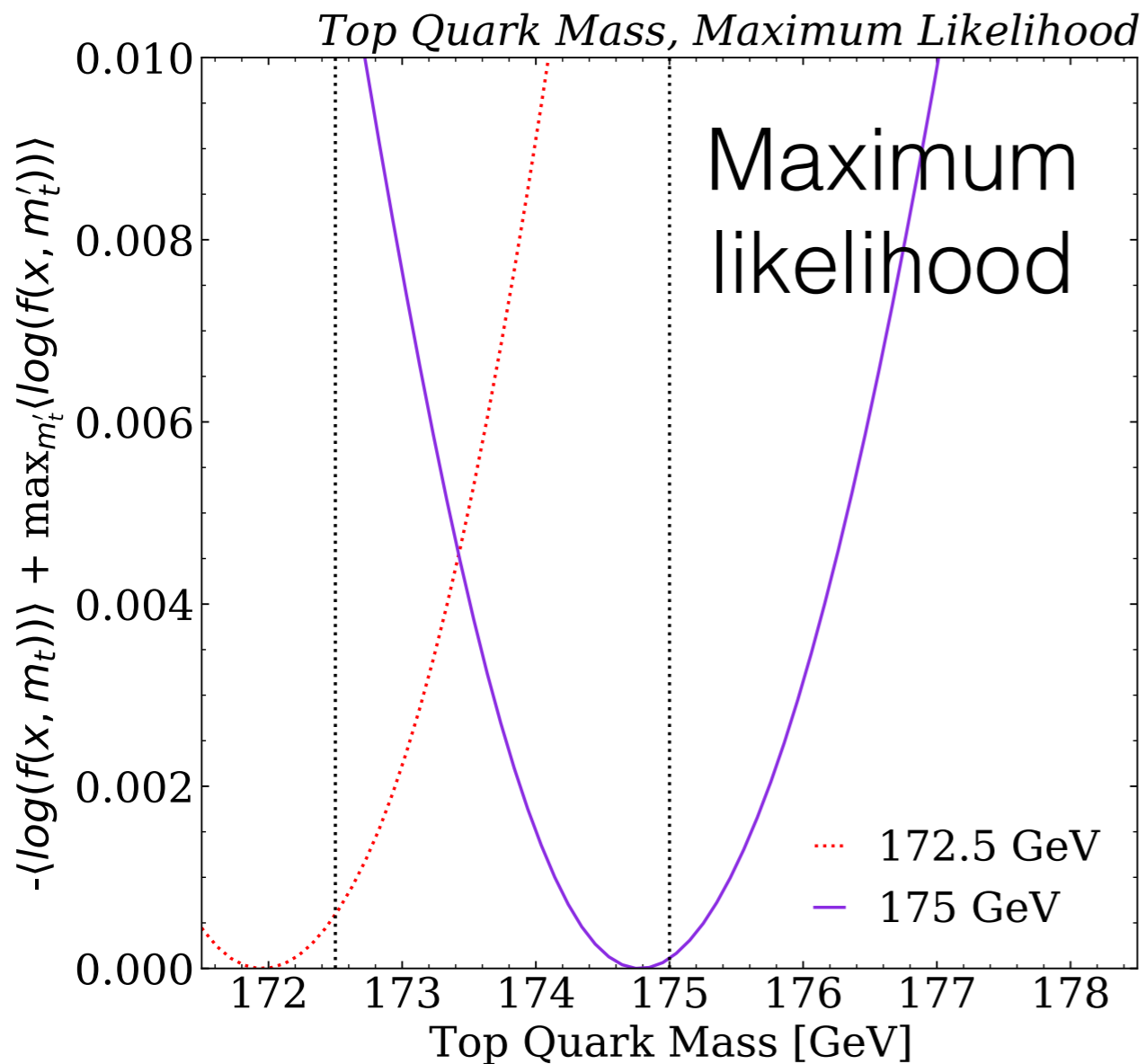
$$g_N(\vec{x}) = \langle \theta | \vec{x} \rangle$$

$$\begin{aligned} g_N(\vec{x}) &= \int d\theta \theta p(\theta | \vec{x}) \\ &= \int d\theta \theta \frac{p(\vec{x} | \theta) p(\theta)}{p(\vec{x})} \\ &= \int d\theta \theta \frac{\frac{p(\vec{x} | \theta)}{p(\vec{x} | \theta_0)} p(\theta)}{\int d\theta' \frac{p(\vec{x} | \theta')}{p(\vec{x} | \theta_0)} p(\theta')} = \int d\theta \theta \frac{f_N(\vec{x}, \theta) p(\theta)}{\int d\theta' f_N(\vec{x}, \theta') p(\theta')} \end{aligned}$$

i.e. this is still secretly per-event classification !

*Gaussian with 100 Instances*





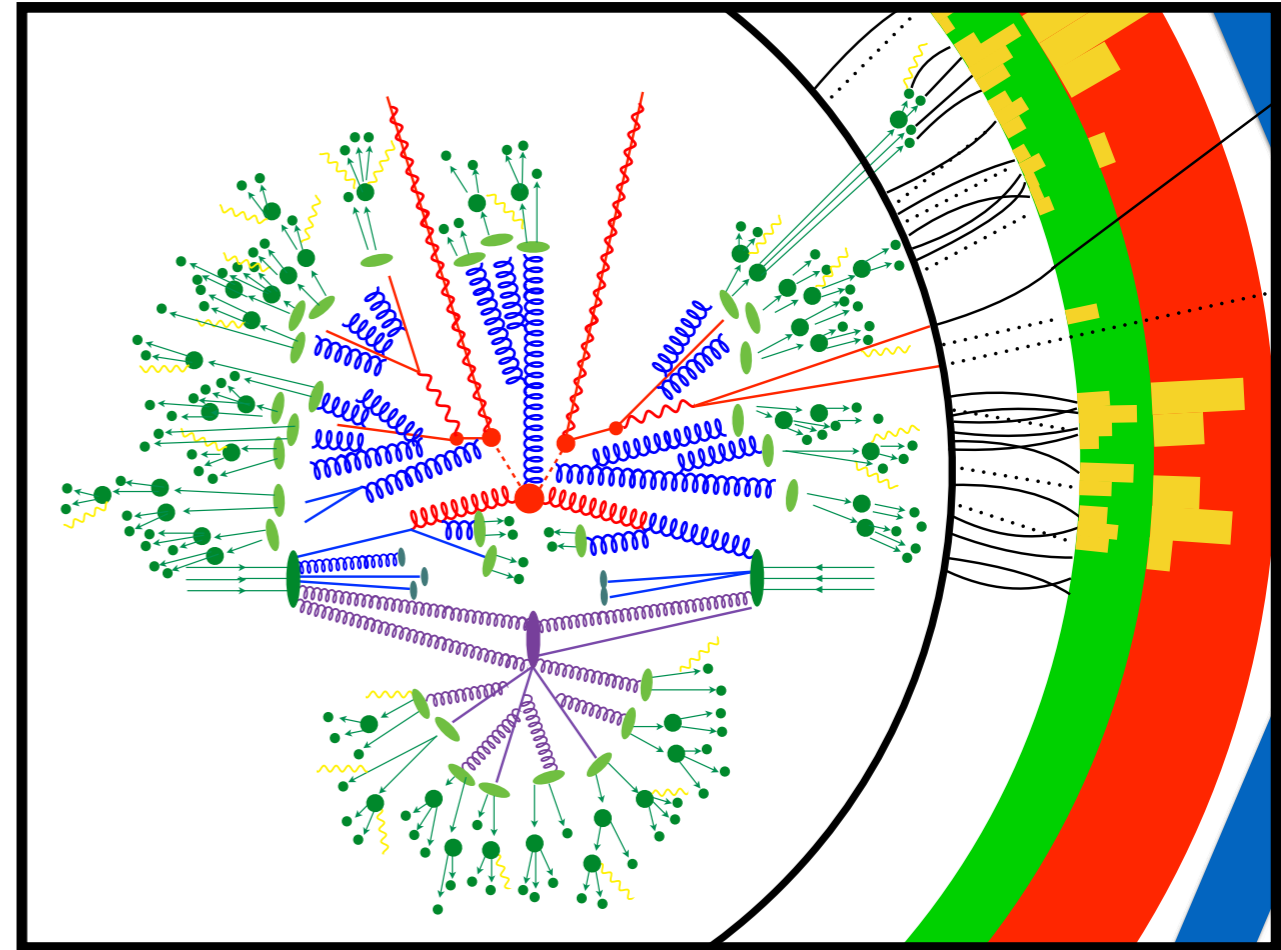
In both cases, the per-instance approaches give the same result as the per-ensemble approach

# Conclusions and outlook

24

Today, I've told you about the interplay between per-event and per-ensemble learning

Formally, these are equivalent. Per-event models are less complex, but there may be practical reasons to prefer one over the other



 [GitHub.com/bnachman/EnsembleLearning](https://github.com/bnachman/EnsembleLearning)

B. Nachman and J. Thaler in PRD 103 (2021) 116013, 2101.07263



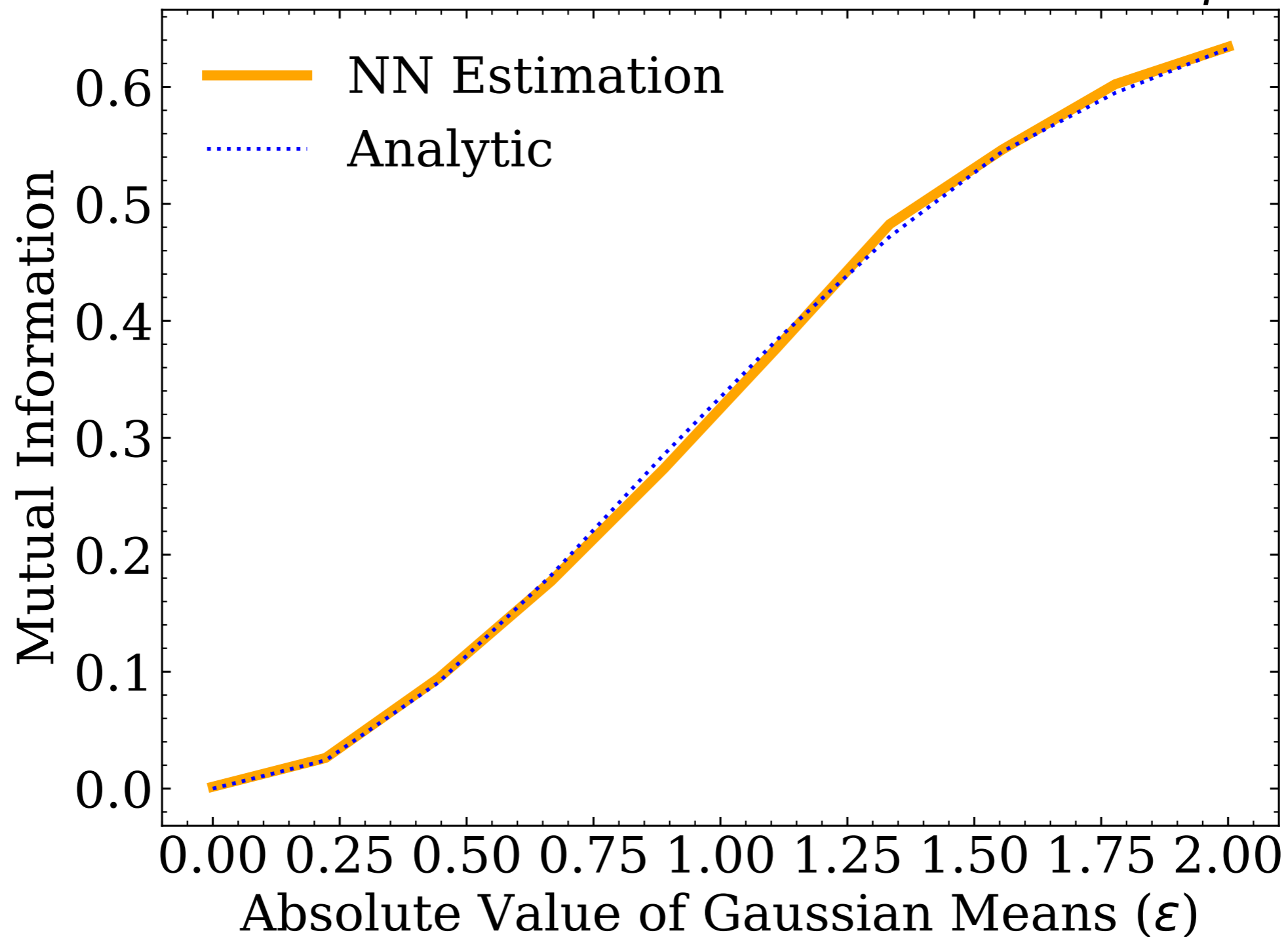
# Backup



# Estimating the Mutual Information

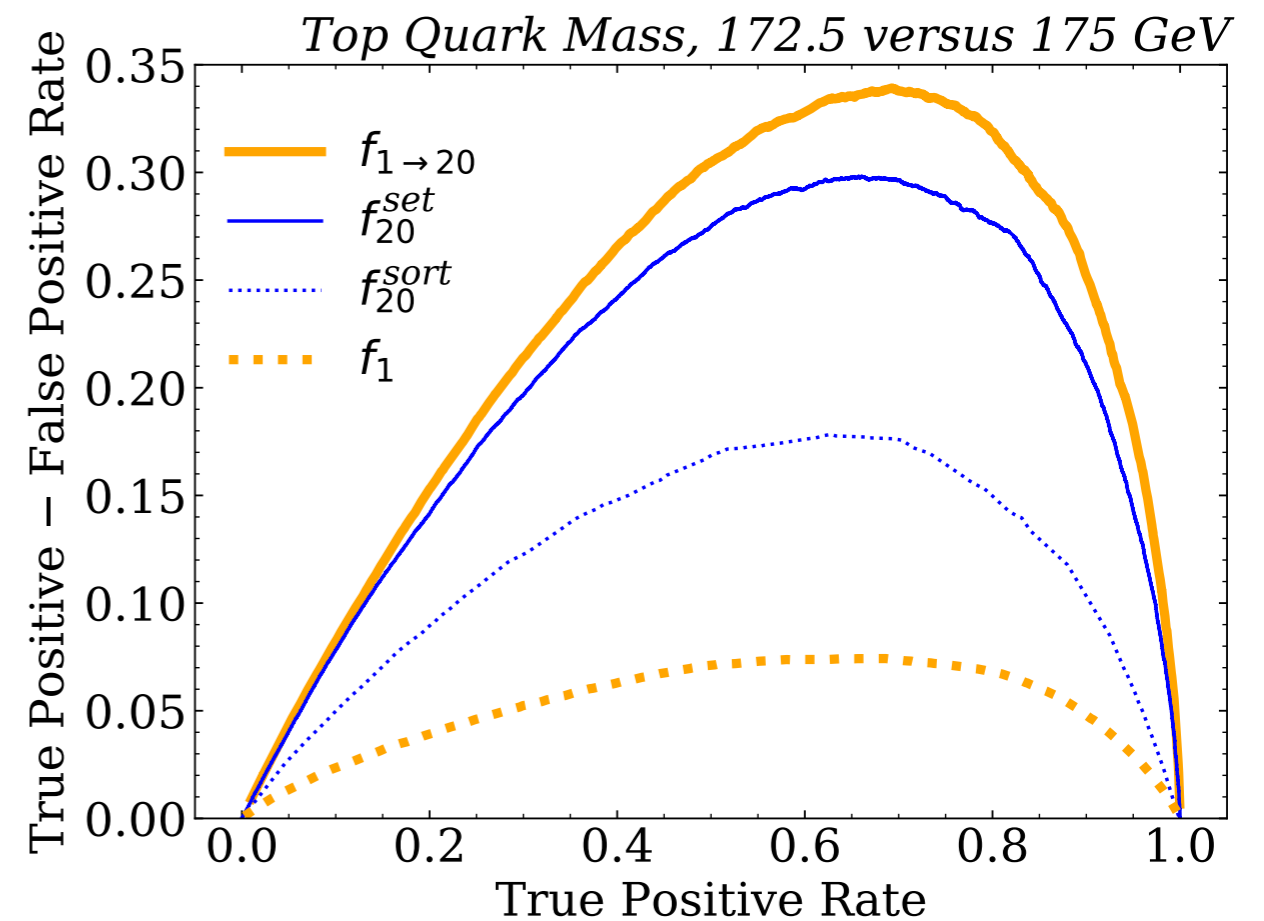
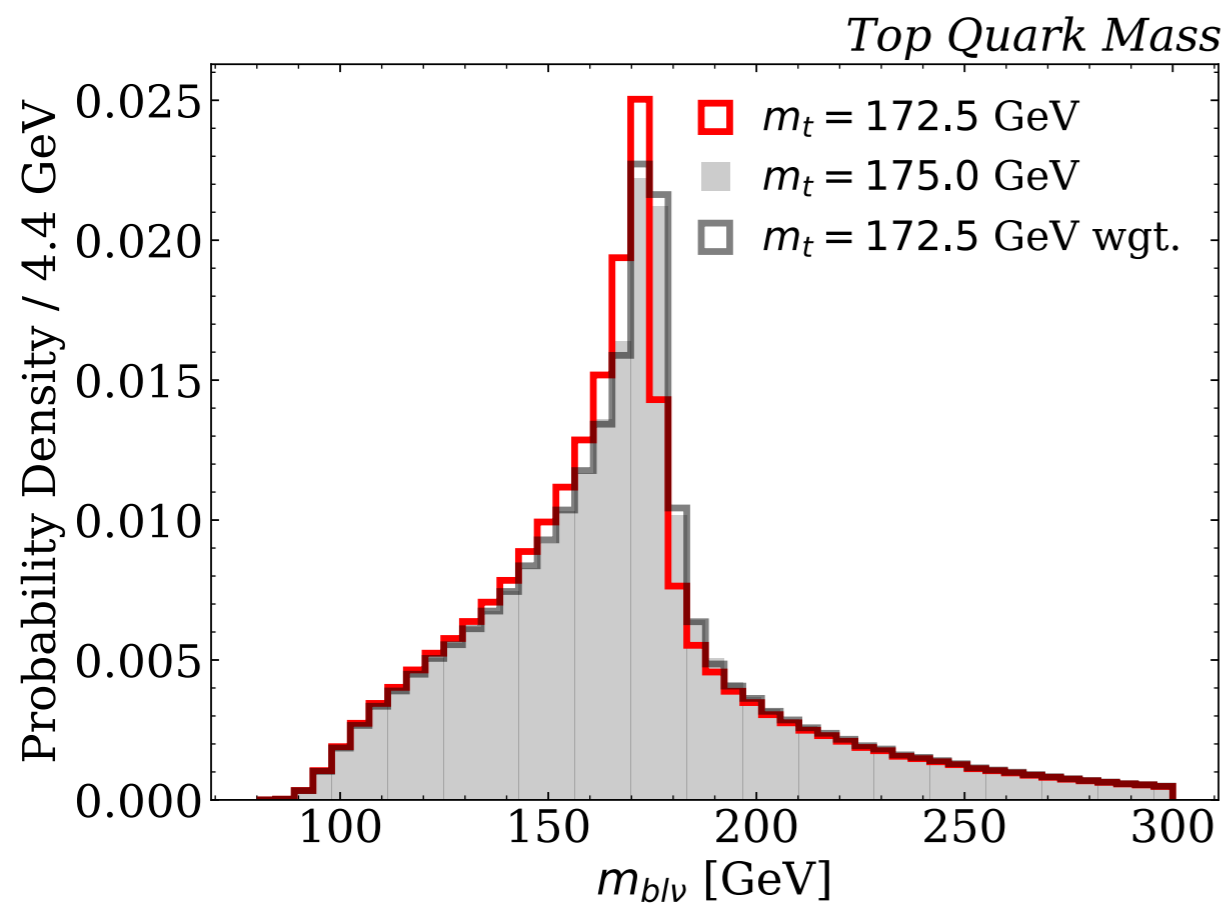
26

*Gaussian Example*



# Examples: Classification [top quarks]

27



# Examples: Classification [top quarks]

28

